

# Prediction of US election with linear model\*

Colin Sihan Yang      Lexun Yu      Siddharth Gowda

October 21, 2024

This paper forecast the winner of the upcoming US presidential election using “poll-of-polls” by building a linear model.

## 1 Introduction

Election result forecasting has become an essential tool for analysts in political science and the public to predict the outcome of democratic process, such as the presidential election in the United States. Traditionally, individual polls have been used as a snapshot of voter sentiment, but they only reflect temporary changes in the performance of contestants, instead of a precise estimation of the election result. As discussed by Pasek (2015) and Blumenthal (2014), the aggregation of multiple polls, or “poll-of-polls,” has become a popular technique to reduce individual survey errors and provide more accurate election forecasts. However, the traditional poll aggregation does not reflect dynamics of an election, especially with real-time changes and the introduction of new data. This creates a gap for a more adaptable model to predict the election result based on both polling data and additional variables, such as historical data and economic indicators.

This paper fills the gap by building a hybrid election forecasting model following the strategies mentioned by Pasek (2015). As Pasek (2015) described in their article, aggregation involves determining which surveys are worth including, as well as selecting, combining and averaging results from multiple polls to reduce individual biases and errors. Prediction modeling adds other data to the model that predicts election outcomes based on current dynamics. Hybrid models like the Bayesian approach incorporates prior beliefs based on historical data or expert knowledge and new evidence like economic updates to dynamically adjust the forecast as the campaign progresses. We incorporate aggregation by filtering the polls on FiveThirtyEight (2024) by numeric grade that indicates pollster’s reliability, prediction that incorporates social

---

\*Code and data are available at: <https://github.com/yulexun/uselection>.

and economic indicators including unemployment rates and abortion rates, and hybrid approaches that leverages Bayesian techniques which combines historical data such as the 2016 election data, allowing for a dynamic prediction of the U.S. presidential election.

The estimand for this research paper is the predicted support percentages for Kamala Harris and Donald Trump. The prediction is based on quantifying various polling factors, including sample size, poll scores, and transparency scores, which are used as predictors.

The results of this model indicate a more stable and accurate forecast compared to traditional aggregation methods alone, [update this ...]

The remainder of this paper is structured as follows: [update this ...]

## 2 Data

### 2.1 Overview

For the data we used in this analysis about the polling result for Kamala Harris and Donald Trump in 2024 USA president election.

- **response variable:** `pct`(`pct`: The percentage of the vote or support that the candidate received in the poll)
- **numeric predictor:**
  - `sample_size`(`sample_size`: The total number of respondents participating in the poll)
  - `timegap`(the time gap between the poll start date and the real election date i.e `timegap = real US election date - poll start date`)
  - `pollscore`(A numeric value representing the score or reliability of the pollster in question)
- **categorical predictor** `state`(The U.S. state where the poll was conducted or focused)
- `methodology`(The method used to conduct the poll)

### 2.2 Measurement

In this dataset, each row represents a polling question that records the variables of interest. Each entry allows us to explore the real-world relationships between polling factors and the support percentage (`pct`) for the candidates Kamala Harris and Donald Trump. This dataset enables an analysis of how various polling characteristics influence the reported support levels for the candidates we are focused.

## 2.3 Clean Data

The data cleaning process involves several steps to ensure the quality and relevance of the polling data. First, we filter the dataset to retain only poll results with a numeric grade of 2.7 or higher, indicating that the polls are considered reliable. Next, we address missing values in the state attribute: polls with NA in the state column are considered national polls.

We then create a new attribute, `days_taken_from_election`, which represents the time gap between the poll's start date and the actual U.S. election date. Additionally, we filter the dataset to include only polls conducted after July 21, 2024, the date when Kamala Harris declared her candidacy. Finally, we remove any remaining rows that contain missing values to ensure a clean dataset.

Table 1: Sample of cleaned US election data

pct	sample_size	pollscore	days_taken_from_election	state	methodology	answer
47.6	4180	-0.8	24	National	Online Ad	Harris
50.7	4180	-0.8	24	National	Online Ad	Trump
0.8	4180	-0.8	24	National	Online Ad	Stein
0.1	4180	-0.8	24	National	Online Ad	Oliver
0.1	4180	-0.8	24	National	Online Ad	West
48.1	4180	-0.8	24	National	Online Ad	Harris

## 2.4 Basic Statistics Summary for Data

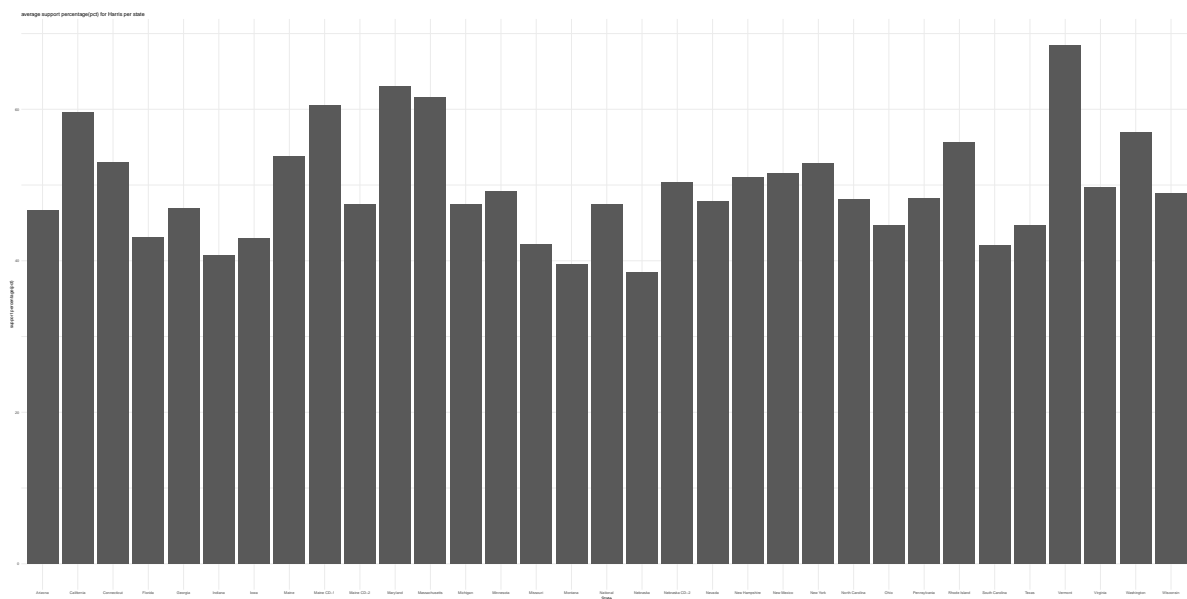
## 3 Model

The goal of our modelling strategy is twofold. Firstly,...

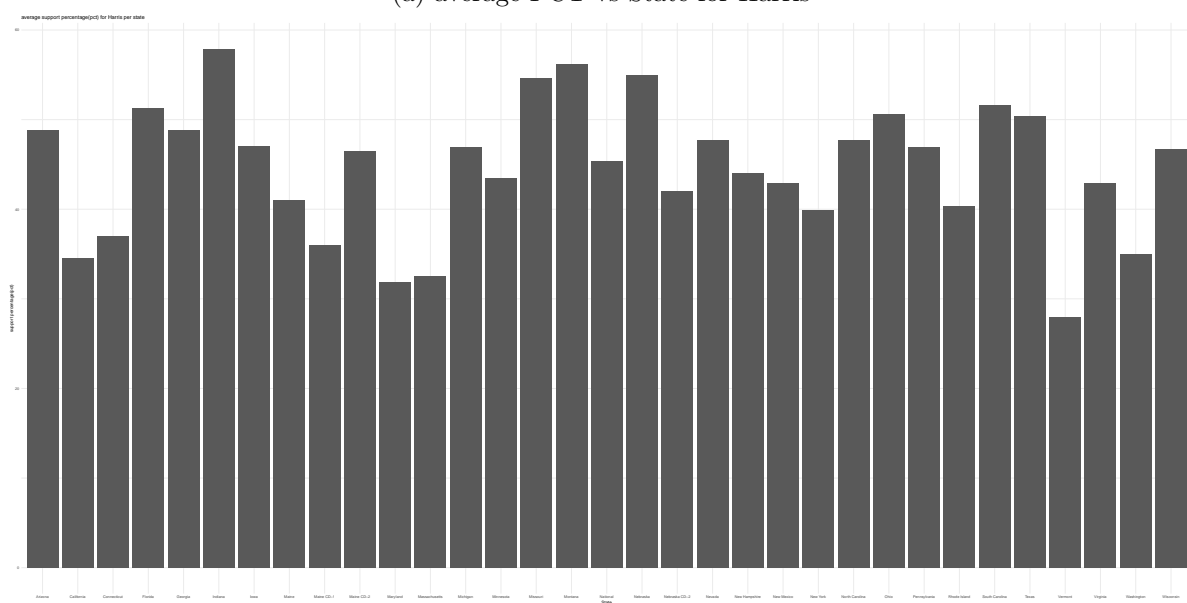
Here we briefly describe the Bayesian analysis model used to investigate... Background details and diagnostics are included in [Appendix B](#).

### 3.1 Model set-up

Define  $y_i$  as the number of seconds that the plane remained aloft. Then  $\beta_i$  is the wing width and  $\gamma_i$  is the wing length, both measured in millimeters.



(a) average PCT vs State for Harris



(b) PCT vs State for Trump

Figure 1: the average PCT vs State for Harris and Trump

$$y_i | \mu_i, \sigma \sim \text{Normal}(\mu_i, \sigma) \tag{1}$$

$$\mu_i = \alpha + \beta_i + \gamma_i \tag{2}$$

$$\alpha \sim \text{Normal}(0, 2.5) \tag{3}$$

$$\beta \sim \text{Normal}(0, 2.5) \tag{4}$$

$$\gamma \sim \text{Normal}(0, 2.5) \tag{5}$$

$$\sigma \sim \text{Exponential}(1) \tag{6}$$

We run the model in R (R Core Team 2023) using the `rstanarm` package of (`rstanarm?`). We use the default priors from `rstanarm`. us

### 3.1.1 Model justification

We expect a positive relationship between the size of the wings and time spent aloft. In particular...

We can use maths by including latex between dollar signs, for instance  $\theta$ .

## 4 Results

Our results are summarized in `?@tbl-modelresults`.

## 5 Discussion

### 5.1 First discussion point

If my paper were 10 pages, then should be be at least 2.5 pages. The discussion is a chance to show off what you know and what you learnt from all this.

### 5.2 Second discussion point

Please don't use these as sub-heading labels - change them to be what your point actually is.

### **5.3 Third discussion point**

### **5.4 Weaknesses and next steps**

Weaknesses and next steps should also be included.

## Appendix

### A Additional data details

### B Model details

#### B.1 Posterior predictive check

In [?@fig-ppcheckandposteriorvsprior-1](#) we implement a posterior predictive check. This shows...

In [?@fig-ppcheckandposteriorvsprior-2](#) we compare the posterior with the prior. This shows...

#### B.2 Diagnostics

[?@fig-stanareyouokay-1](#) is a trace plot. It shows... This suggests...

[?@fig-stanareyouokay-2](#) is a Rhat plot. It shows... This suggests...

### C FiftyEight Licenses

FiftyEight's [data sets](#) are used and modified by us under the [Creative Commons Attribution 4.0 International License](#).

## References

- Blumenthal, Mark. 2014. “Polls, Forecasts, and Aggregators.” *PS: Political Science and Politics* 47 (2): 297–300. <http://www.jstor.org/stable/43284537>.
- FiveThirtyEight. 2024. “Our Data.” *FiveThirtyEight*. <https://data.fivethirtyeight.com>.
- Pasek, Josh. 2015. “THE POLLS–REVIEW: PREDICTING ELECTIONS: CONSIDERING TOOLS TO POOL THE POLLS.” *The Public Opinion Quarterly* 79 (2): 594–619. <http://www.jstor.org/stable/24546379>.
- R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.